

オンライン学習(15章)

杉山将・本多淳也

sugi@k.u-tokyo.ac.jp, jhonda@k.u-tokyo.ac.jp

<http://www.ms.k.u-tokyo.ac.jp>

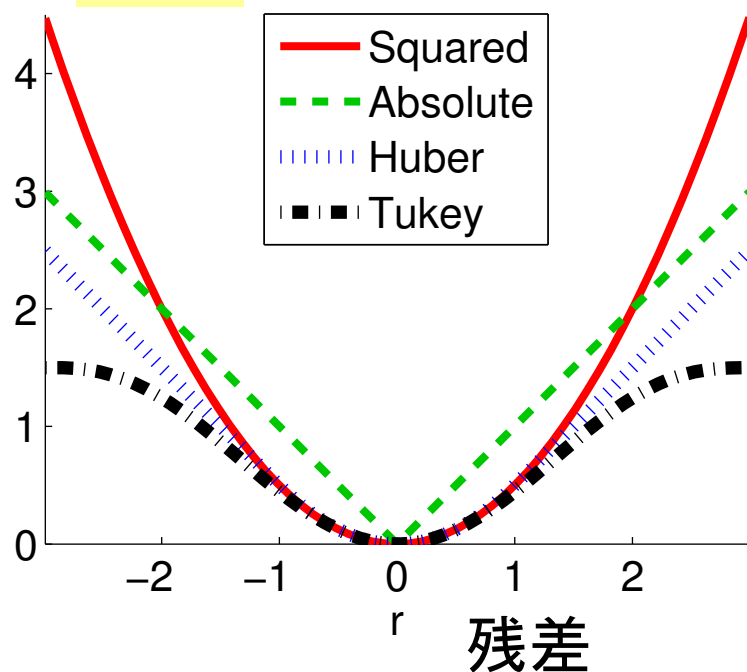
一括学習

2

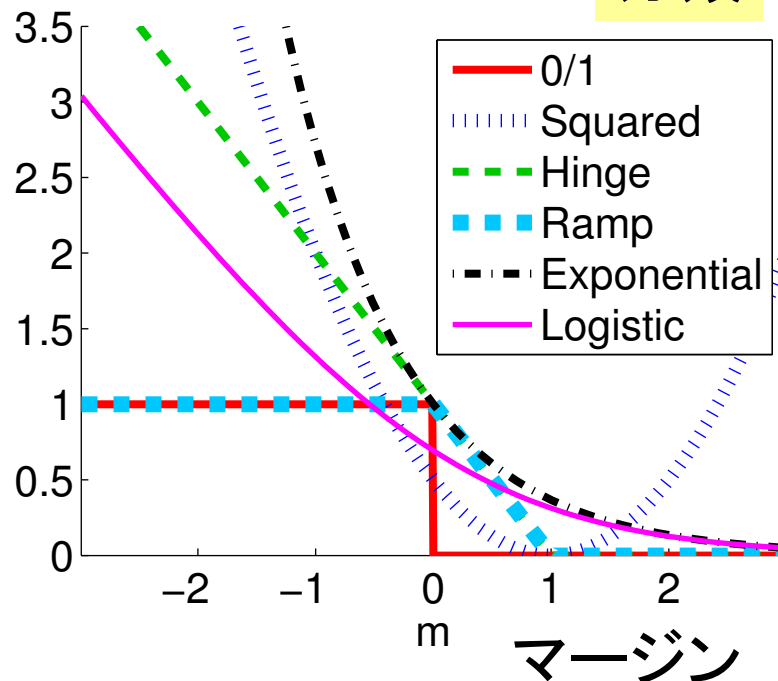
- 全ての訓練標本を用いてパラメータを学習
- 標本が多いとき、一括で学習するのは困難

$$\min_{\theta} \sum_{i=1}^n \text{loss} \left(f_{\theta}(\mathbf{x}_i), y_i \right) \quad \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

回帰



分類



- データを一つ一つ(あるいはいくつかつ)学習していく
- 一回一回の解の更新は非常に簡単に行える
- データを生成する確率分布が変化するような場面にも適応できる

講義の流れ

1. 確率的勾配法
2. 受動攻撃学習
3. 適応正則化学習



確率的勾配法

5

$$\min_{\theta} \sum_{i=1}^n \text{loss}(f_{\theta}(\mathbf{x}_i), y_i)$$

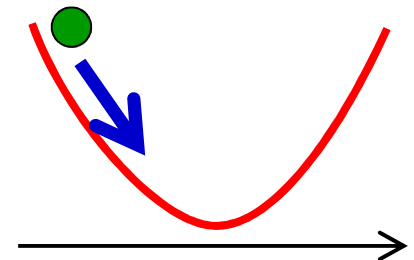
$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x})$$

1. パラメータ θ を適当に初期化
2. 標本 (\mathbf{x}, y) をランダムに選び, 勾配を少し降下

$$\theta \leftarrow \theta - \varepsilon \nabla \text{loss}(f_{\theta}(\mathbf{x}), y)$$

ステップ幅
 $\varepsilon > 0$

3. 収束するまで 2. を繰り返す
- 実用上は, 一つでなく複数の標本(ミニバッチ)を用いることが多い



実行例

■ ガウスカーネルモデル

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{j=1}^n \theta_j K(\boldsymbol{x}, \boldsymbol{x}_j)$$

$$K(\boldsymbol{x}, \boldsymbol{c}) = \exp \left(-\frac{\|\boldsymbol{x} - \boldsymbol{c}\|^2}{2h^2} \right)$$

に対する最小二乗法

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i=1}^n \left(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - y_i \right)^2$$

の確率的勾配法

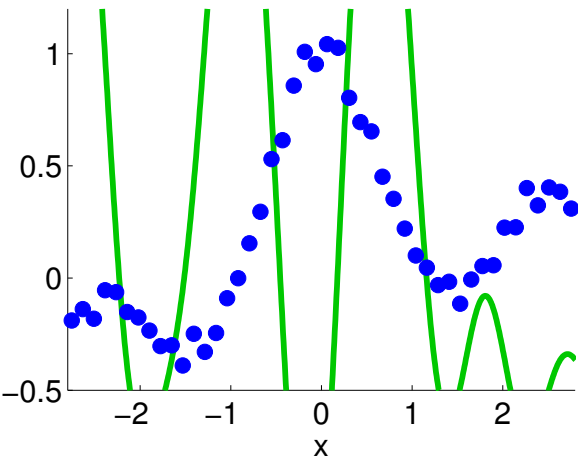
■ 勾配は $k \left(k^{\top} \boldsymbol{\theta} - y \right)$

$$\boldsymbol{k} = (K(\boldsymbol{x}, \boldsymbol{x}_1), \dots, K(\boldsymbol{x}, \boldsymbol{x}_n))^{\top}$$

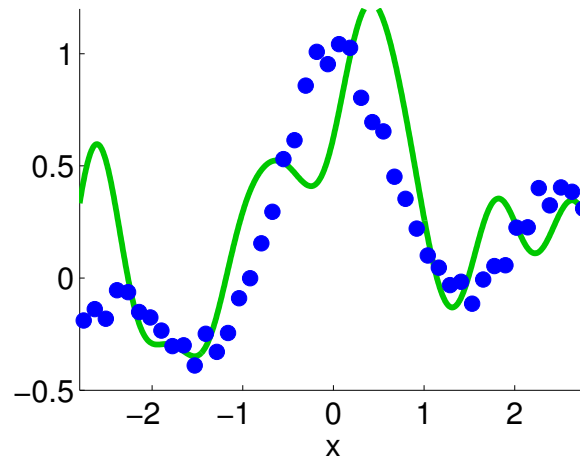
実行例(続き)

7

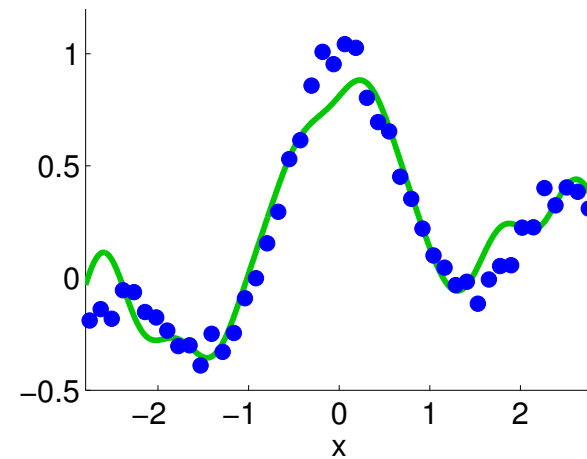
初期値



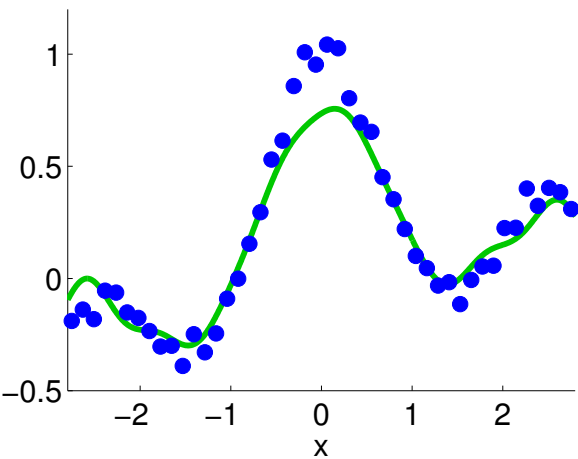
50反復後



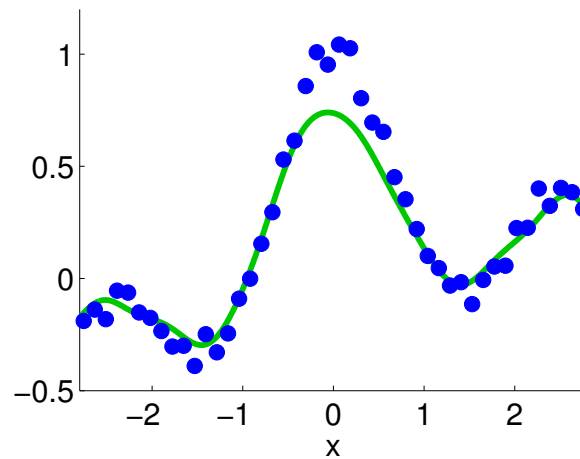
100反復後



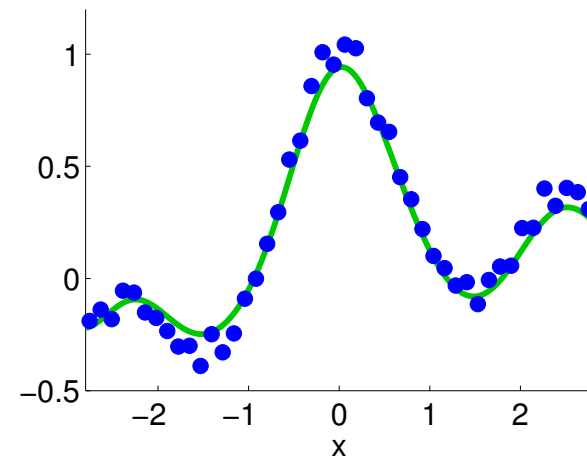
150反復後



200反復後



11556反復後



実装例

8

```
clear all; rand('state',0); randn('state',0);
n=50; N=1000; x=linspace(-3,3,n)';
X=linspace(-3,3,N)';
pix=pi*x; y=sin(pix)./(pix)+0.1*x+0.05*randn(n,1);

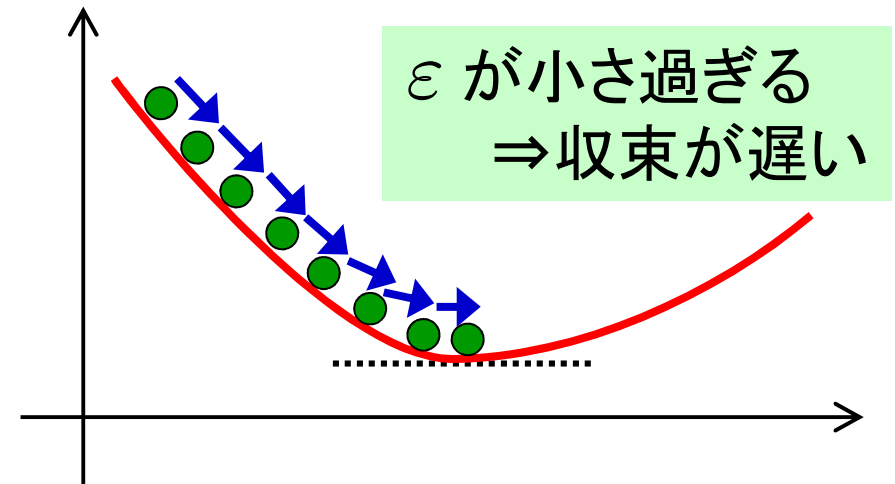
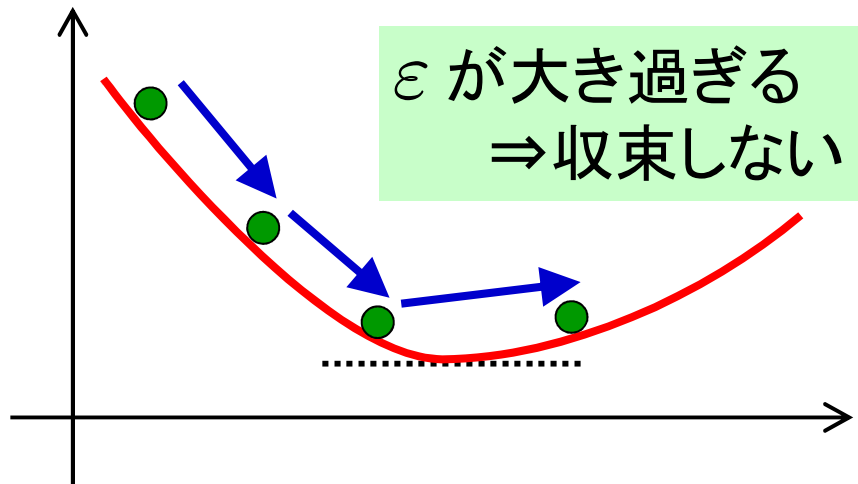
hh=2*0.3^2; t0=randn(n,1); e=0.1;
for o=1:n*1000
    i=ceil(rand*n); ki=exp(-(x-x(i)).^2/hh);
    t=t0-e*ki*(ki'*t0-y(i));
    if norm(t-t0)<0.000001, break, end
    t0=t;
end
K=exp(-(repmat(X.^2,1,n)+repmat(x.^2',N,1)-2*X*x')/hh);
F=K*t;

figure(1); clf; hold on; axis([-2.8 2.8 -0.5 1.2]);
plot(X,F,'g-'); plot(x,y,'bo');
```


確率的勾配法の問題点

9

■ ステップ幅 ε の設定が難しい



- 最小二乗回帰では、一気に谷底まで勾配降下することが可能(停留解が解析的に求まる)
- しかし、そのような急激な学習を行うと、過去の学習で得られた知識が損なわれる恐れがある

講義の流れ

1. 確率的勾配法
2. 受動攻撃学習
3. 適応正則化学習

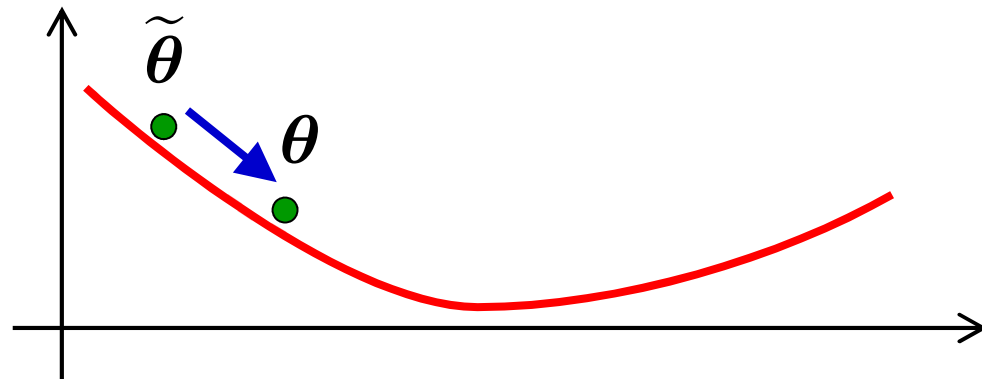


受動攻撃(passive-aggressive)学習¹¹

- 現在の解 $\tilde{\theta}$ からの変化量を抑制する

$$\min_{\theta} \underbrace{\text{loss}(f_{\theta}(x), y)}_{\text{訓練出力に対する適合の良さ}} + \underbrace{\gamma \|\theta - \tilde{\theta}\|^2}_{\text{パラメータの変化量を抑制}}$$

受動係数
 $\gamma > 0$



基底関数
 $\phi(x)$

- 以下, 線形モデルを考える: $f_{\theta}(x) = \theta^{\top} \phi(x)$

確率的勾配法との比較

- 二乗損失に対する受動攻撃回帰の解の更新式:

$$\theta \leftarrow \theta - \frac{\theta^\top \phi(x) - y}{\|\phi(x)\|^2 + \gamma} \phi(x) \quad \gamma > 0$$

導出は後ほど

$$\text{loss}(f_\theta(x), y) = (f_\theta(x) - y)^2$$

- 確率的勾配法の解の更新式:

$$\theta \leftarrow \theta - \varepsilon \phi(x) (\theta^\top \phi(x) - y)$$

- 受動攻撃回帰では, ステップ幅をデータ x に合わせて適応的に調整:

$$\varepsilon = \frac{1}{\|\phi(x)\|^2 + \gamma}$$

二乗ヒンジ損失

■ 二乗ヒンジ損失に対する受動攻撃分類

$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta)$$

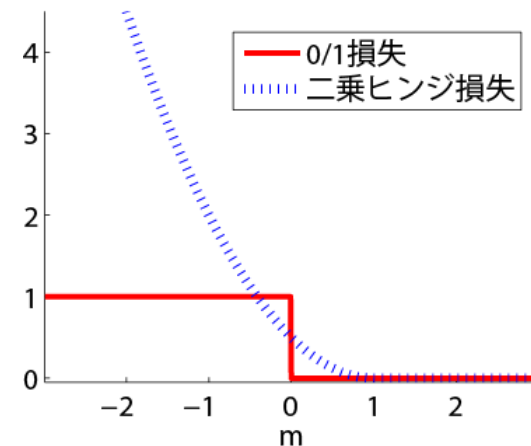
$$\gamma > 0$$

$$J(\theta) = \frac{1}{2} \left(\max \left(0, 1 - \theta^{\top} \phi(x) y \right) \right)^2 + \frac{\gamma}{2} \|\theta - \tilde{\theta}\|^2$$

の解は次式で与えられる

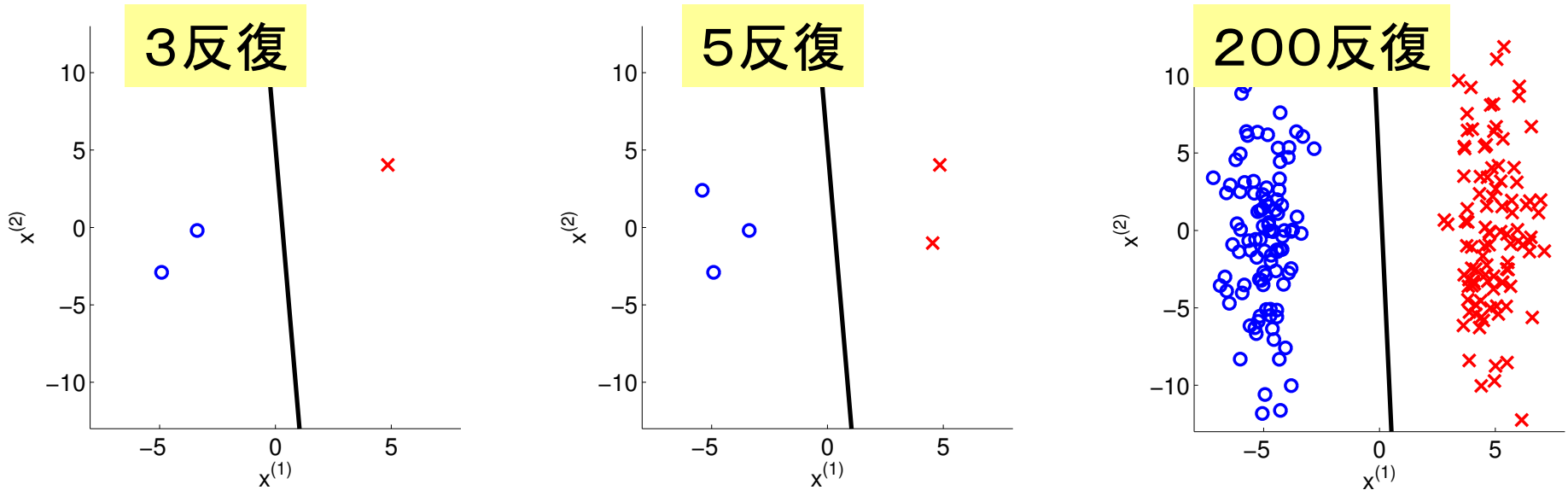
$$\theta \leftarrow \theta + \frac{y \max \left(0, 1 - y \theta^{\top} \phi(x) \right)}{\|\phi(x)\|^2 + \gamma} \phi(x)$$

導出は後ほど

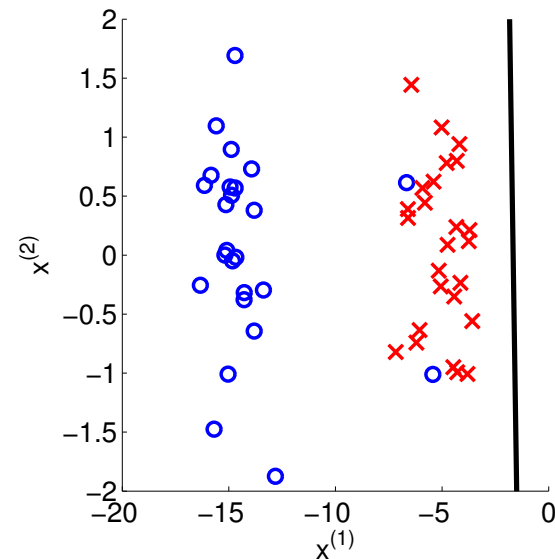


実行例

■ うまく分類できている



■ ただし、異常値に弱い

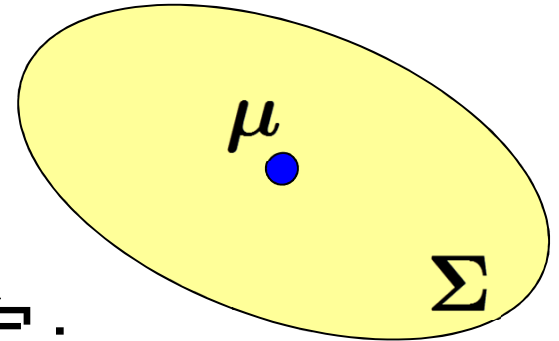


講義の流れ

1. 確率的勾配法
2. 受動攻撃学習
3. 適応正則化学習



- パラメータを推定するかわりに
パラメータの確率分布を推定
(ベイズ推定とは異なる)



- 簡単のため, **ガウス分布**を仮定:

$$\theta \sim N(\mu, \Sigma)$$

- パラメータの変化量を分布間の**カルバック・ライブラー距離**で測る:

$$\text{KL}[N(\mu, \Sigma) \| N(\tilde{\mu}, \tilde{\Sigma})]$$

$$\text{KL}(p \| \tilde{p}) = \int p(\theta) \log \frac{p(\theta)}{\tilde{p}(\theta)} d\theta$$

- 分散共分散行列の適応正則化: $\phi(x)^\top \Sigma \phi(x)$

■ パラメータの分布がガウス分布

$$\theta \sim N(\mu, \Sigma)$$

にしたがうとき, 点 x での予測値 $f_{\theta}(x) = \theta^{\top} \phi(x)$ の分散は

$$\text{Var}[f_{\theta}(x)] = \phi(x)^{\top} \Sigma \phi(x)$$

で与えられることを示せ

■ ヒント: 分散の定義

$$\text{Var}[X] = E[(X - E[X])(X - E[X])]$$

■ パラメータの学習規準:

$$\min_{\mu, \Sigma} \left\{ \text{loss} \left(f_{\mu}(\mathbf{x}), y \right) \right.$$

訓練出力に対する
適合の良さ

$$+ 2\gamma \text{KL} \left[N(\mu, \Sigma) \parallel N(\tilde{\mu}, \tilde{\Sigma}) \right]$$

パラメータの
変化量を抑制

$$\left. + \phi(\mathbf{x})^{\top} \Sigma \phi(\mathbf{x}) \right\}$$

適応正則化:
予測値の分散を抑制

$$\gamma > 0$$

$$\text{KL}(p \parallel \tilde{p}) = \int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{\tilde{p}(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

■ ガウス分布間のカルバック・ライブラー距離が

$$\text{KL} \left[N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \parallel N(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \right]$$

$$= \frac{1}{2} \left\{ \log \frac{\det(\tilde{\boldsymbol{\Sigma}})}{\det(\boldsymbol{\Sigma})} + \text{tr}(\tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}) - d \right\}$$

で与えられることを示せ.

$d: x$ の次元数

■ ガウス分布 $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ の確率密度関数 $p(\boldsymbol{\theta})$ は

$$p(\boldsymbol{\theta}) = (2\pi)^{-d/2} \det(\boldsymbol{\Sigma})^{-1/2} \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) \right)$$

■ $(\boldsymbol{\theta} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) = \text{tr} \left(\boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) (\boldsymbol{\theta} - \boldsymbol{\mu})^\top \right)$

■ 分散共分散行列: $\int p(\boldsymbol{\theta}) (\boldsymbol{\theta} - \boldsymbol{\mu}) (\boldsymbol{\theta} - \boldsymbol{\mu})^\top d\boldsymbol{\theta} = \boldsymbol{\Sigma}$

■ 期待値ベクトル: $\int p(\boldsymbol{\theta}) \boldsymbol{\theta} d\boldsymbol{\theta} = \boldsymbol{\mu}$

■ 以下の式を示せ:

$$\text{KL} \left[N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \parallel N(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \right]$$

$$\text{KL}(p \parallel \tilde{p}) = \int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{\tilde{p}(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

$$= \frac{1}{2} \left\{ \log \frac{\det(\tilde{\boldsymbol{\Sigma}})}{\det(\boldsymbol{\Sigma})} + \text{tr}(\tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}) - d \right\}$$

$$p(\boldsymbol{\theta}) = (2\pi)^{-d/2} \det(\boldsymbol{\Sigma})^{-1/2} \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) \right)$$

$$(\boldsymbol{\theta} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) = \text{tr} \left(\boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) (\boldsymbol{\theta} - \boldsymbol{\mu})^\top \right)$$

$$\int p(\boldsymbol{\theta}) \boldsymbol{\theta} d\boldsymbol{\theta} = \boldsymbol{\mu}$$

$$\int p(\boldsymbol{\theta}) (\boldsymbol{\theta} - \boldsymbol{\mu}) (\boldsymbol{\theta} - \boldsymbol{\mu})^\top d\boldsymbol{\theta} = \boldsymbol{\Sigma}$$

受動攻撃学習との関係

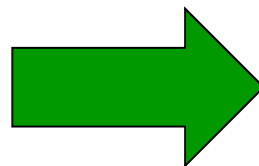
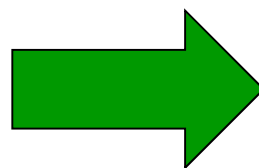
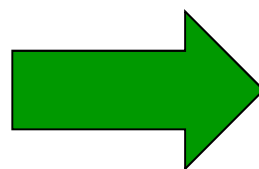
- Σ を単位行列に固定すれば,
 μ の学習規準は受動攻撃学習と一致:

適応正則化学習

$$\text{loss}(f_{\mu}(x), y)$$

$$\text{KL}[N(\mu, \Sigma) \parallel N(\tilde{\mu}, \tilde{\Sigma})]$$

$$\phi(x)^{\top} \Sigma \phi(x)$$



受動攻撃学習

$$\text{loss}(f_{\theta}(x), y)$$

$$\|\theta - \tilde{\theta}\|^2$$

定数

解の求め方: 共分散行列

25

$$(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} J(\mu, \Sigma)$$

$$J(\mu, \Sigma) = \operatorname{loss}\left(f_{\mu}(x), y\right) + \phi(x)^{\top} \Sigma \phi(x)$$

$$+ \gamma \left\{ \log \frac{\det(\tilde{\Sigma})}{\det(\Sigma)} + \operatorname{tr}(\tilde{\Sigma}^{-1} \Sigma) + (\mu - \tilde{\mu})^{\top} \tilde{\Sigma}^{-1} (\mu - \tilde{\mu}) - d \right\}$$

- **数学演習**: 適応正則化学習の Σ の解が次式で与えられることを示せ

$$\hat{\Sigma} = \tilde{\Sigma} - \frac{\tilde{\Sigma} \phi(x) \phi(x)^{\top} \tilde{\Sigma}}{\phi(x)^{\top} \tilde{\Sigma} \phi(x) + \gamma}$$

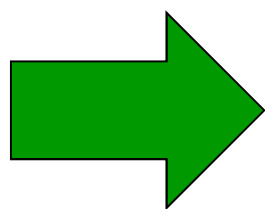
ヒント

■ 行列での微分の公式を利用する:

- $\frac{\partial}{\partial \Sigma} \phi(x)^\top \Sigma \phi(x) = \phi(x) \phi(x)^\top$
- $\frac{\partial}{\partial \Sigma} \log \det(\Sigma) = \Sigma^{-1}$
- $\frac{\partial}{\partial \Sigma} \text{tr} \left(\tilde{\Sigma}^{-1} \Sigma \right) = \tilde{\Sigma}^{-1}$

■ 逆行列の公式を利用する:

$$\hat{\Sigma}^{-1} = \tilde{\Sigma}^{-1} + \frac{1}{\gamma} \phi(x) \phi(x)^\top$$



$$\hat{\Sigma} = \tilde{\Sigma} - \frac{\tilde{\Sigma} \phi(x) \phi(x)^\top \tilde{\Sigma}}{\phi(x)^\top \tilde{\Sigma} \phi(x) + \gamma}$$

解の求め方: 共分散行列(再掲) ²⁷

$$J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \text{loss}\left(f_{\boldsymbol{\mu}}(\boldsymbol{x}), y\right) + \boldsymbol{\phi}(\boldsymbol{x})^{\top} \boldsymbol{\Sigma} \boldsymbol{\phi}(\boldsymbol{x}) + \gamma \left\{ \log \frac{\det(\tilde{\boldsymbol{\Sigma}})}{\det(\boldsymbol{\Sigma})} + \text{tr}(\tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}})^{\top} \tilde{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}) - d \right\}$$

■ 数学演習: 以下の式を示せ

$$\hat{\boldsymbol{\Sigma}} = \tilde{\boldsymbol{\Sigma}} - \frac{\tilde{\boldsymbol{\Sigma}} \boldsymbol{\phi}(\boldsymbol{x}) \boldsymbol{\phi}(\boldsymbol{x})^{\top} \tilde{\boldsymbol{\Sigma}}}{\boldsymbol{\phi}(\boldsymbol{x})^{\top} \tilde{\boldsymbol{\Sigma}} \boldsymbol{\phi}(\boldsymbol{x}) + \gamma}$$

$$(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) = \underset{\boldsymbol{\mu}, \boldsymbol{\Sigma}}{\text{argmin}} J(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \boldsymbol{\phi}(\boldsymbol{x})^{\top} \boldsymbol{\Sigma} \boldsymbol{\phi}(\boldsymbol{x}) = \boldsymbol{\phi}(\boldsymbol{x}) \boldsymbol{\phi}(\boldsymbol{x})^{\top}$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \log \det(\boldsymbol{\Sigma}) = \boldsymbol{\Sigma}^{-1}$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \text{tr} \left(\tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma} \right) = \tilde{\boldsymbol{\Sigma}}^{-1}$$

$$\left(\tilde{\boldsymbol{\Sigma}}^{-1} + \frac{1}{\gamma} \boldsymbol{\phi}(\boldsymbol{x}) \boldsymbol{\phi}(\boldsymbol{x})^{\top} \right)^{-1} = \tilde{\boldsymbol{\Sigma}} - \frac{\tilde{\boldsymbol{\Sigma}} \boldsymbol{\phi}(\boldsymbol{x}) \boldsymbol{\phi}(\boldsymbol{x})^{\top} \tilde{\boldsymbol{\Sigma}}}{\boldsymbol{\phi}(\boldsymbol{x})^{\top} \tilde{\boldsymbol{\Sigma}} \boldsymbol{\phi}(\boldsymbol{x}) + \gamma}$$

解の求め方: 期待値

- **数学演習**: 二乗損失に対する適応正則化回帰の μ の解は次式で与えられることを示せ

$$\hat{\mu} = \tilde{\mu} - \frac{\phi(x)^\top \tilde{\mu} - y}{\phi(x)^\top \tilde{\Sigma} \phi(x) + \gamma} \tilde{\Sigma} \phi(x)$$

$$(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} J(\mu, \Sigma)$$

$$J(\mu, \Sigma) = \left(\mu^\top \phi(x) - y \right)^2 + \phi(x)^\top \Sigma \phi(x)$$

$$+ \gamma \left\{ \log \frac{\det(\tilde{\Sigma})}{\det(\Sigma)} + \operatorname{tr}(\tilde{\Sigma}^{-1} \Sigma) + (\mu - \tilde{\mu})^\top \tilde{\Sigma}^{-1} (\mu - \tilde{\mu}) - d \right\}$$

- **ヒント**: $\mu^\top \phi(x) = \phi(x)^\top \mu$ および逆行列の公式

$$\left(\tilde{\Sigma}^{-1} + \frac{1}{\gamma} \phi(x) \phi(x)^\top \right)^{-1} = \tilde{\Sigma} - \frac{\tilde{\Sigma} \phi(x) \phi(x)^\top \tilde{\Sigma}}{\phi(x)^\top \tilde{\Sigma} \phi(x) + \gamma}$$

適応正則化回帰

■ パラメータの更新式:

$$\mu \longleftarrow \mu - (\mu^\top \phi(x) - y) \Sigma \phi(x) / \beta$$

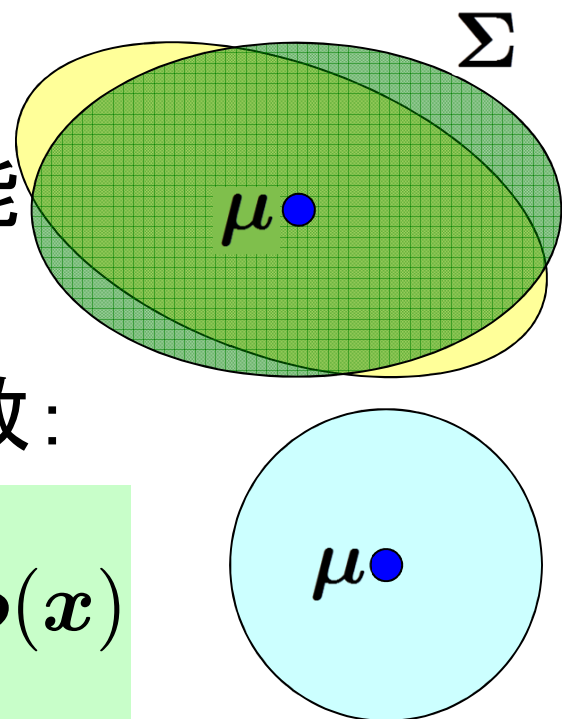
$$\Sigma \longleftarrow \Sigma - \Sigma \phi(x) \phi(x)^\top \Sigma / \beta$$

$$\beta = \phi(x)^\top \Sigma \phi(x) + \gamma$$

■ Σ を**対角行列**で近似すれば,
計算時間とメモリ使用量を削減可能

■ Σ を**単位行列**に固定すれば,
 μ の更新式は受動攻撃学習と一致:

$$\theta \longleftarrow \theta - \frac{\theta^\top \phi(x) - y}{\|\phi(x)\|^2 + \gamma} \phi(x)$$



二乗ヒンジ損失に対する 期待値の解

$$(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} J(\mu, \Sigma)$$

$$J(\mu, \Sigma) = \left(\max \left(0, 1 - \mu^\top \phi(x)y \right) \right)^2 + \phi(x)^\top \Sigma \phi(x)$$

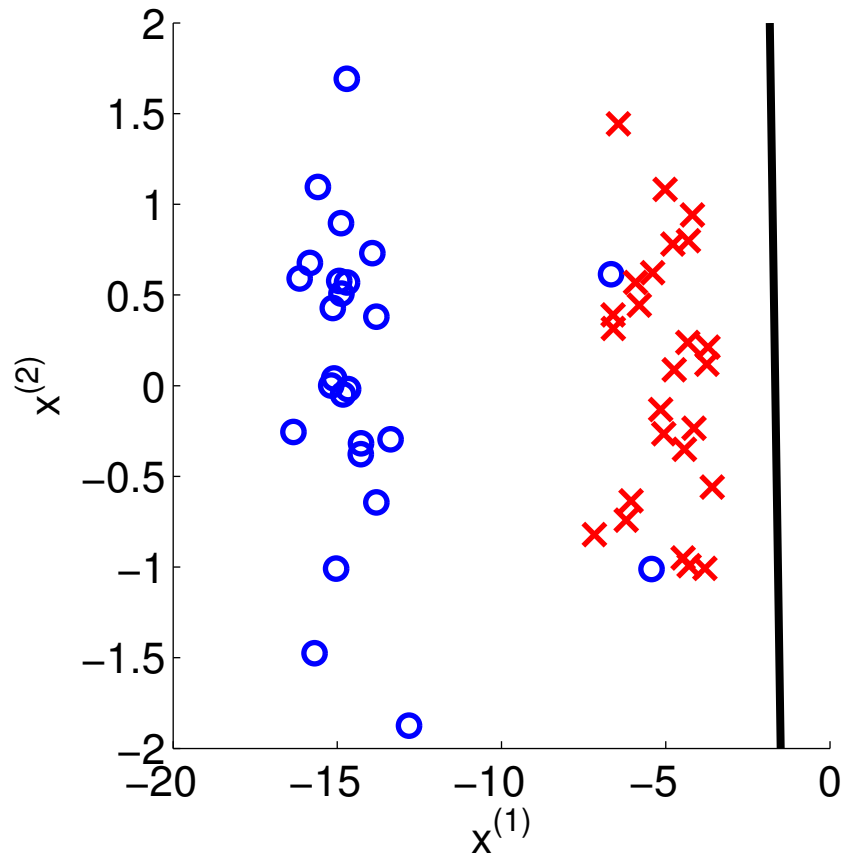
$$+ \gamma \left\{ \log \frac{\det(\tilde{\Sigma})}{\det(\Sigma)} + \operatorname{tr}(\tilde{\Sigma}^{-1} \Sigma) + (\mu - \tilde{\mu})^\top \tilde{\Sigma}^{-1} (\mu - \tilde{\mu}) - d \right\}$$

- 二乗ヒンジ損失に対する適応正則化分類の μ の解は次式で与えられる:

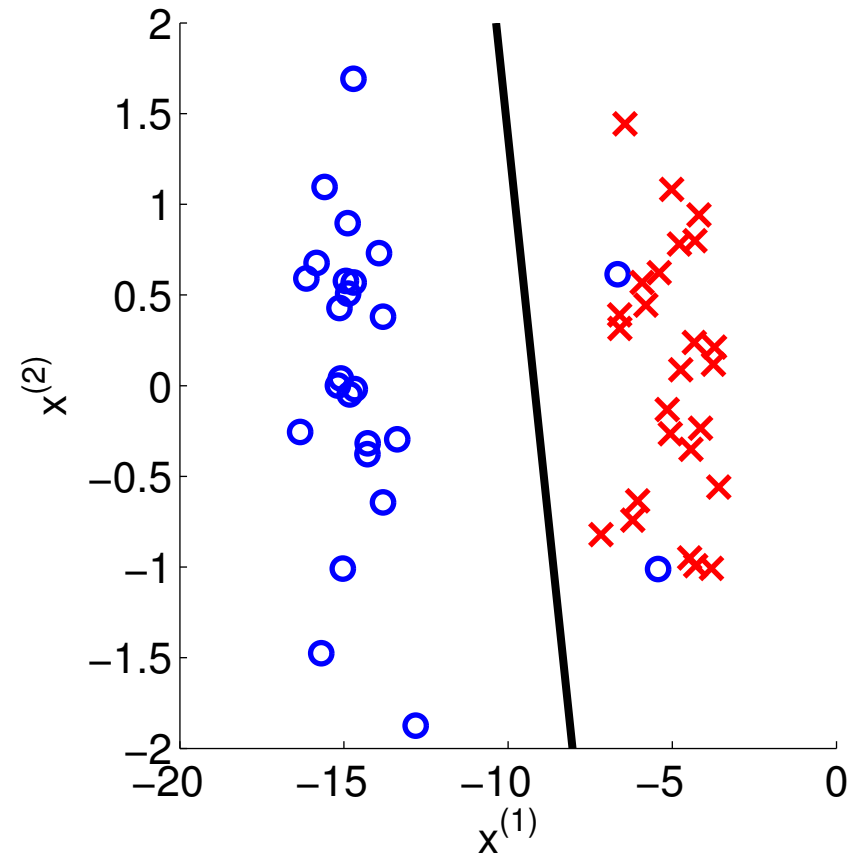
$$\hat{\mu} = \tilde{\mu} + \frac{y \max(0, 1 - \tilde{\mu}^\top \phi(x)y)}{\phi(x)^\top \tilde{\Sigma} \phi(x) + \gamma} \tilde{\Sigma} \phi(x)$$

証明は
宿題

受動攻撃分類



適応正則化分類

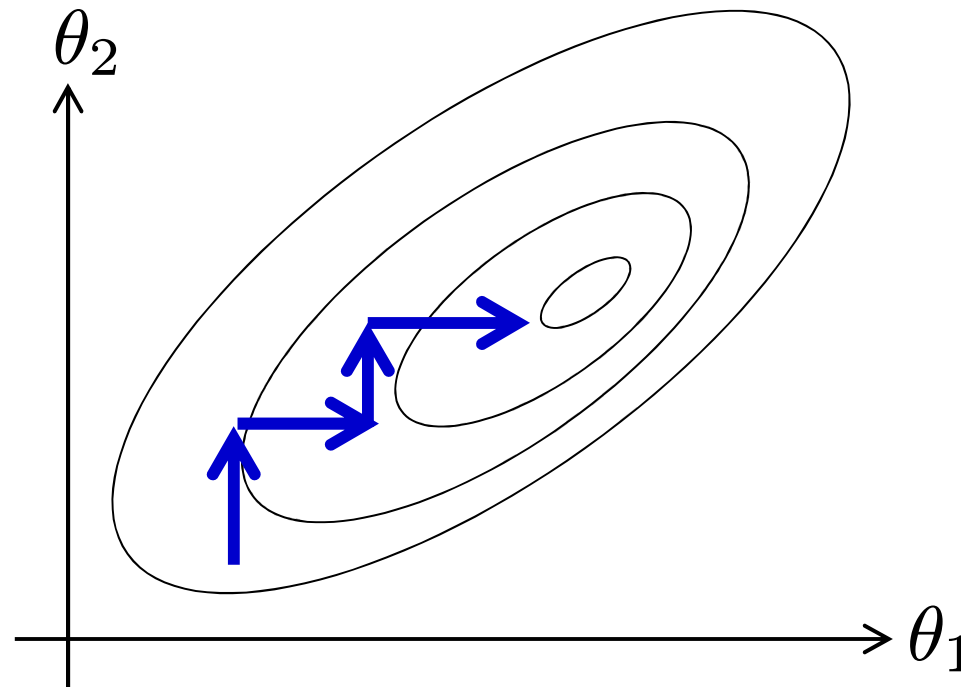


■ 異常値に対してロバストな解が得られている

オンライン学習:まとめ

36

- オンライン学習では, 一部の標本のみを使うことにより, 勾配の計算を高速化
- パラメータを一部分ごとに学習する**座標降下法**も実用上有用



講義の流れ



1. 確率的勾配法
2. 受動攻撃学習
3. 適応正則化学習

- オンライン学習 : データを一つ一つ学習していく
 - 確率的勾配法 : 学習係数の決定が困難
 - 受動攻撃学習 : パラメータの変化量を抑制
 - 適応正則化学習 : パラメータの分布を考える

次回の予告

- 半教師付き学習(16章)
- 転移学習(18章)



$$(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} J(\mu, \Sigma)$$

$$J(\mu, \Sigma) = \left(\max \left(0, 1 - \mu^\top \phi(x)y \right) \right)^2 + \phi(x)^\top \Sigma \phi(x)$$

$$+ \gamma \left\{ \log \frac{\det(\tilde{\Sigma})}{\det(\Sigma)} + \operatorname{tr}(\tilde{\Sigma}^{-1} \Sigma) + (\mu - \tilde{\mu})^\top \tilde{\Sigma}^{-1} (\mu - \tilde{\mu}) - d \right\}$$

- 二乗ヒンジ損失に対する適応正則化分類の μ の解は次式で与えられることを示せ

$$\hat{\mu} = \tilde{\mu} + \frac{y \max(0, 1 - \tilde{\mu}^\top \phi(x)y)}{\phi(x)^\top \tilde{\Sigma} \phi(x) + \gamma} \tilde{\Sigma} \phi(x)$$

■ データ

```
clear all; rand('state',0); randn('state',0); n=50;  
x=[randn(1,n/2)-15 randn(1,n/2)-5; randn(1,n)]';  
y=[ones(n/2,1); -ones(n/2,1)]; x(1:2,1)=x(1:2,1)+10;  
x(:,3)=1; p=randperm(n); x=x(p,:); y=y(p);
```

を用いて, 二乗ヒンジ損失に
基づく適応正則化分類を
線形モデル

$$\begin{aligned} f_{\theta}(x) &= f_{\theta}(x^{(1)}, x^{(2)}) \\ &= (x^{(1)} \ x^{(2)} \ 1)\theta \end{aligned}$$

に対して実装せよ

